

ECE421 Introduction to Machine Learning

Assignment #3

Unsupervised Learning and Probabilistic Models

Yan, Xuanming
xuanming.yan@mail.utoronto.ca
1001175382
Contribution: 50%

Wang, Chu Qing
eliza.wang@mail.utoronto.ca
1001303573
Contribution: 50%

Mar 2019

1 K-means

1.1

We had implemented the distance function by using pairwise squared distance between X and MU , and the helper function "assign_data" returns a list of indexes for every x to its closest K-mean. We chose 200 iterations to train the data, since we could see obvious decay in loss.

First, we randomly initialize k centers using Gaussian distribution, with $stddev = 1$. $\mu^k \in \mathbb{R}^D$. Then Assign each point $k \in \{1, \dots, K\}$ to nearest center: $C^{(t)}(n) = \operatorname{argmin}_{k=1}^K \|\mathbf{x}_n - \mu_k\|_2^2$. Finally, update μ_i , let it becomes new centroid.

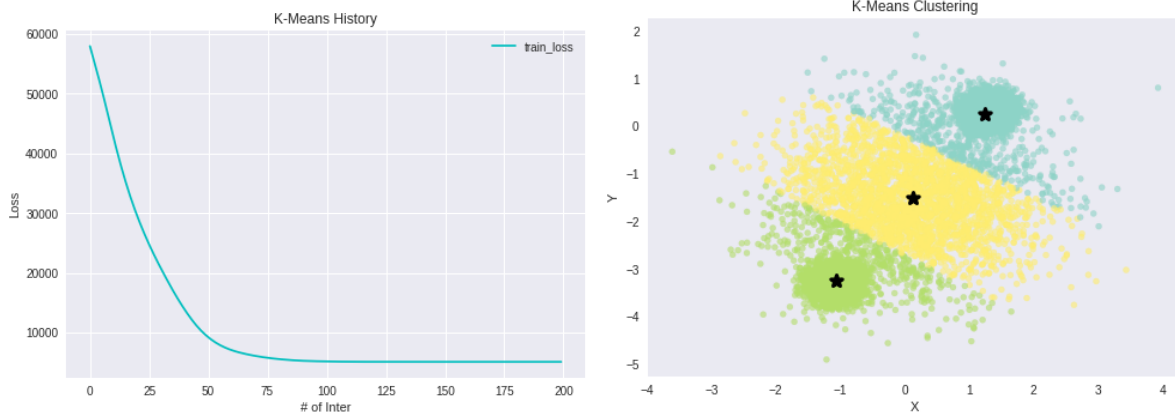


Figure 1: K-means when $K = 3$. Losses vs Iterations

```
1 import tensorflow as tf
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import helper as hlp
5
6
7 def assign_data(X, MU):
8
9     dists = distanceFunc(X, MU)
10    min_dist = tf.argmin(dists, 1)
11    return min_dist
12
13 # Distance function for K-means
14 def distanceFunc(X, MU):
15     # Inputs
16     # X: is an NxD matrix (N observations and D dimensions)
```

```

17 # MU: is an KxD matrix (K means and D dimensions)
18 # Outputs
19 # pair_dist: is the pairwise distance matrix (NxK)
20 #y = np.power((X[:, np.newaxis] - MU), 2)
21 #newY = tf.convert_to_tensor(y, dtype=tf.float32)
22 #output = np.sum(y, axis=2)
23 newX = tf.expand_dims(X,0)
24 newMU = tf.expand_dims(MU, 1)
25 dis = tf.reduce_sum(tf.square(tf.subtract(newX,newMU)),2)
26 output = tf.transpose(dis)
27 return output
28
29 # Loading data
30 data = np.load('data2D.npy')
31 # data = np.load('data100D.npy')
32 [num_pts, dim] = np.shape(data)
33
34 # For Validation set
35 is_valid = False
36 if is_valid:
37     valid_batch = int(num_pts / 3.0)
38     np.random.seed(45689)
39     rnd_idx = np.arange(num_pts)
40     np.random.shuffle(rnd_idx)
41     val_data = data[rnd_idx[:valid_batch]]
42     data = data[rnd_idx[valid_batch:]]
43
44 #graph = tf.Graph()
45 #with graph.as_default():
46 loss_history = np.empty(shape=[0], dtype=float)
47 val_loss_history = []
48 K = 5
49 N = num_pts
50 D = dim
51 inter = 201
52
53 X = tf.placeholder("float", shape=[None,D])
54 MU_init = tf.truncated_normal([K,D], stddev=0.25)
55 MU = tf.Variable(MU_init)
56
57 distance = distanceFunc(X,MU)
58 #centroid = np.power(distance,2)
59 #loss = np.sum(np.amin(distance, axis=1))
60 loss = tf.reduce_sum(tf.reduce_min(distance, axis = 1))
61 optimizer = tf.train.AdamOptimizer(
62     learning_rate=0.05, beta1=0.9, beta2=0.99, epsilon=1e-5).minimize(loss)
63
64 init_g = tf.global_variables_initializer()
65 sess = tf.Session()
66 sess.run(init_g)
67
68 for step in range(inter):
69     cenVal,lossVal,_ = sess.run([MU,loss,optimizer], feed_dict={X:data})
70     val_cenVal,val_lossVal,_ = sess.run([MU,loss,optimizer], feed_dict={X:val_data})
71     loss_history = np.append(loss_history,lossVal)
72     val_loss_history = np.append(val_loss_history, val_lossVal)
73     if step%10 ==0:
74         print("iteration:", step, "loss", lossVal)
75     clustering = sess.run(assign_data(X, MU), feed_dict={X: data, MU:cenVal})
76     print("Validation loss", val_lossVal)
77     percentages = np.zeros(K)
78     for i in range(K):
79         percentages[i] = np.sum(np.equal(i, clustering))*100.0/len(clustering)
80         print("class:", i, "percentage:", percentages[i])
81 plt.figure(1)
82 plt.plot(range(len(loss_history)),loss_history,c="c", label="train_loss")
83 plt.plot(range(len(val_loss_history)),val_loss_history,c="b", label="validation_loss")
84 plt.legend(loc = "best")
85 plt.title('K-Means History')
86 plt.xlabel('# of Inter')
87 plt.ylabel('Loss')
88 plt.show()
89
90 k = len(cenVal)

```

```

91 plt.scatter(data[:, 0], data[:, 1], c=clustering,
92             cmap=plt.get_cmap('Set3'), s=25, alpha=0.6)
93 plt.scatter(cenVal[:, 0], cenVal[:, 1], marker='*', c="black",
94             cmap=plt.get_cmap('Set1'), s=50, linewidths=1)
95 plt.title('K-Means Clustering')
96 plt.xlabel('X')
97 plt.ylabel('Y')
98 plt.grid()
99 plt.show()

```

Listing 1: K-means using Adam optimizer

1.2

By looking at table 1 below, the percentage of data in each clusters are balanced up to $k = 3$, from $K = 4$ and onward, there exists unbalanced percentage data in each clusters, some clusters has much more data accumulated than others, therefore $k = 3$ would be the best choice.

	Class 1	Class 2	Class 3	Class 4	Class 5
K = 1	100%	0	0	0	0
K = 2	49.53%	50.47%	0	0	0
K = 3	23.81%	38.06%	38.13%	0	0
K = 4	12.03%	13.53%	37.13%	37.31%	0
K = 5	11.07%	7.63%	37.04%	7.55%	36.71%

Table 1: Percentage of the data points belonging to each of the K clusters.

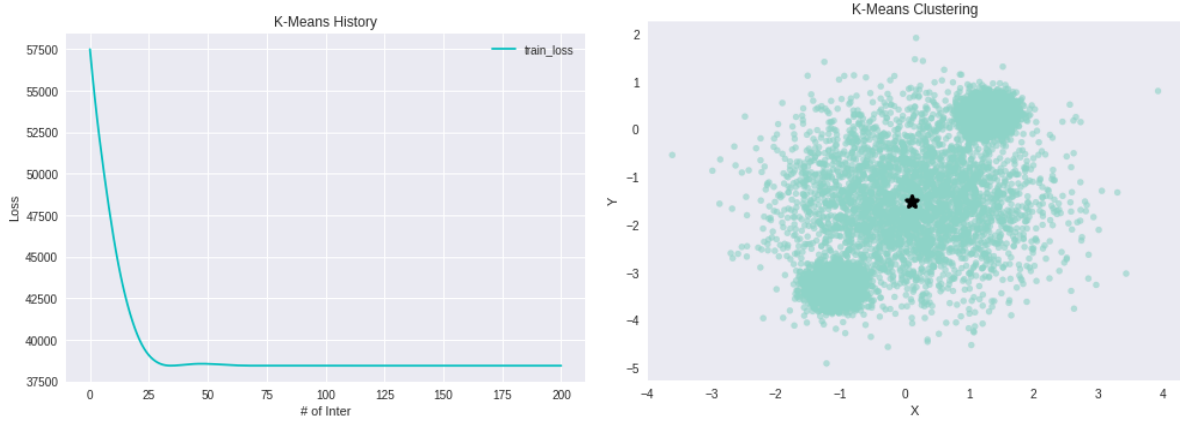


Figure 2: K-means when $K = 1$. Losses vs Iterations

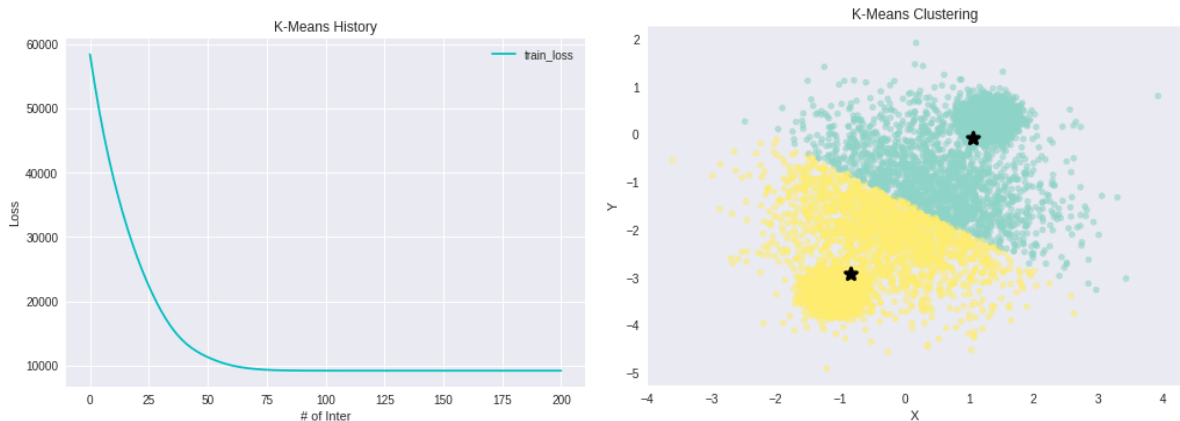


Figure 3: K-means when $K = 2$. Losses vs Iterations

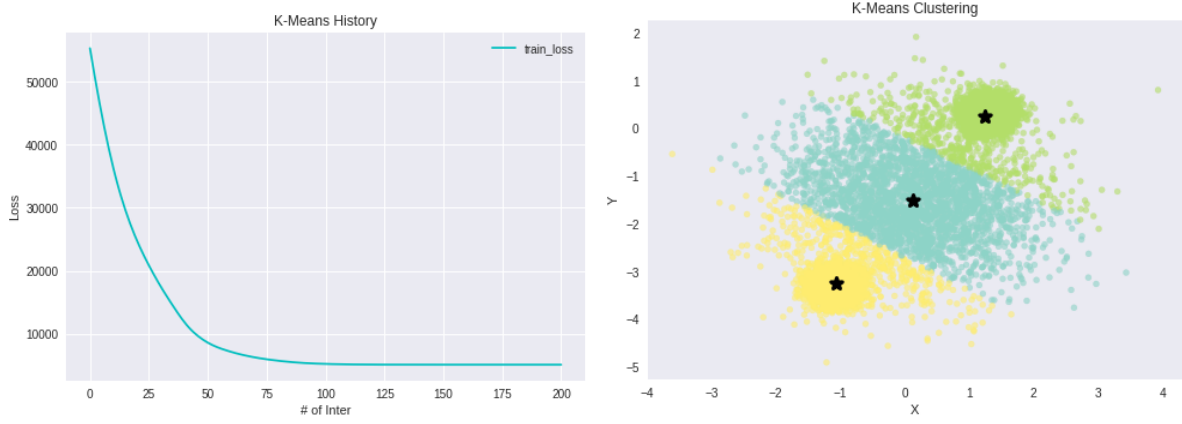


Figure 4: K-means when $K = 3$. Losses vs Iterations



Figure 5: K-means when $K = 4$. Losses vs Iterations

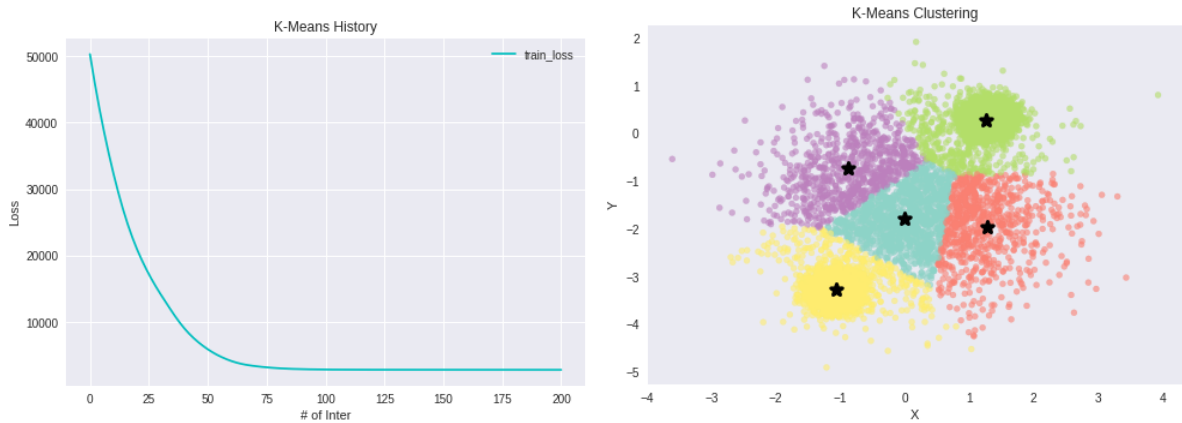


Figure 6: K-means when $K = 5$. Losses vs Iterations

1.3

We cannot draw conclusion by only looking the the loss for different cluster size, since more clusters will always gives us smaller loss value, if number of $k = \text{number of data}$, we can even get loss equals to zero. Therefore, instead of looking at the actual loss values, we can look at the trend of change in loss values as shown in figure 7. We notice that there is not much change in the slope after $k = 3$, therefore, $k = 3$ should be the best choice.

Number of clusters	1	2	3	4	5
Validation loss	12857.85	2959.22	1613.18	1053.37	885.79

Table 2: Loss of validation data with respect to different cluster number.



Figure 7: Change in validation loss with respect to cluster numbers

2 Mixture of Gaussians

2.1 The Gaussian cluster mode

2.1.1 Log Probability density function implementation

Randomly initialize K centers, we have:

$$\begin{aligned}
 \mathcal{N}(x|\mu_k, \sigma_k^2) &= P(x|\mu_k, \sigma_k^2) \\
 \log \mathcal{N}(x|\mu_k, \sigma_k^2) &= \log P(x|\mu_k, \sigma_k^2) \\
 &= \log\left(\frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{(x - \mu_k)^2}{2\sigma_k^2}\right)\right) \\
 &= \log\left(\frac{1}{\sqrt{2\pi}\sigma_k}\right) + \left(-\frac{(x - \mu_k)^2}{2\sigma_k^2}\right)
 \end{aligned}$$

Python implementation of above equation show below.

```

1 def log_GaussPDF(X, mu, sigma):
2     # Inputs
3     # X: N X D
4     # mu: K X D
5     # sigma: K X 1
6     # log_pi: K X 1
7
8     # Outputs:
9     # log_Gaussian PDF N X K
10    dim = tf.to_float(tf.rank(X))
11    xmu = distanceFunc(X, mu)
12    xmuSqu = tf.multiply(xmu, xmu)
13    sigma = tf.squeeze(sigma)
14    coef = tf.log(2 * np.pi * sigma)
15    exp = xmu / (2 * sigma)
16    mul = -0.5 * dim * coef
17    PDF = mul - exp
18    return PDF

```

Listing 2: Log probability density

2.1.2 Log Probability of the cluster variable z given data vector x

Basic idea of this part is to invert the probability found in part 2.1.1. Apply Bayes's rule, we can simply get:

$$\begin{aligned} P(z = k|x) &= \frac{P(x|z = k)P(z = k)}{P(x)} \\ &= \frac{P(x|z = k)P(z = k)}{\sum_i P(x|z = i)P(z = i)} \end{aligned}$$

Where in above equation, $\sum_{i=1}^{i=K} P(x|z = i)P(z = i)$ gives all possible scenarios of clustering. Therefore,

$$\log P(z = k|x) = \log P(x|z = k) + \log P(z = k) - \log \sum_i \exp(\log P(x|z = i) + \log P(z = i))$$

```

1 def log_posterior(log_PDF, log_pi):
2     # Input
3     # log_PDF: log Gaussian PDF N X K
4     # log_pi: K X 1
5
6     # Outputs
7     # log_post: N X K
8     log_pi = tf.squeeze(log_pi)
9     log_prob = tf.add(log_pi, log_PDF)
10    log_sum = hlp.reduce_logsumexp(log_prob + log_pi, keep_dims=True)
11    output = log_prob - log_sum
12    return output

```

Listing 3: Log posterior clustering probability

2.2 Learning the MoG

2.2.1

The reasonable initial value for parameters were shown in Listing 4 below. We had enforced the constraint for π by using a softmax function, and enforced σ in the range of 0 to inf by doing exp to σ . All parameters were initialized by sampling from the standard normal distribution, with the standard deviation of 0.01, and table 3 shows the reason for this selection.

```

1 learning_rate = 0.003
2 s_stddev=0.05
3 X = tf.placeholder("float", [None, D], "X")
4 mu = tf.Variable(tf.random_normal([K, D], stddev = s_stddev))
5 sigma = tf.Variable(tf.random_normal([K, 1], stddev = s_stddev))
6 sigma = tf.exp(sigma)
7 log_PDF = log_GaussPDF(X, mu, sigma)
8
9 initial_pi = tf.Variable(tf.random_normal([K, 1], stddev = s_stddev))
10 log_pi = tf.squeeze(hlp.logsoftmax(initial_pi))

```

Listing 4: Initialization of parameters.

σ	Beginning Loss	Best Loss
0.01	49016.14	17132.400
1	52995.58	17132.543
10	359978000	82524.844

Table 3: Choice of standard deviation value when initializing parameters

With these parameters and selection of $k = 3$, $iteration = 1000$. The result of trained test data had shown below.

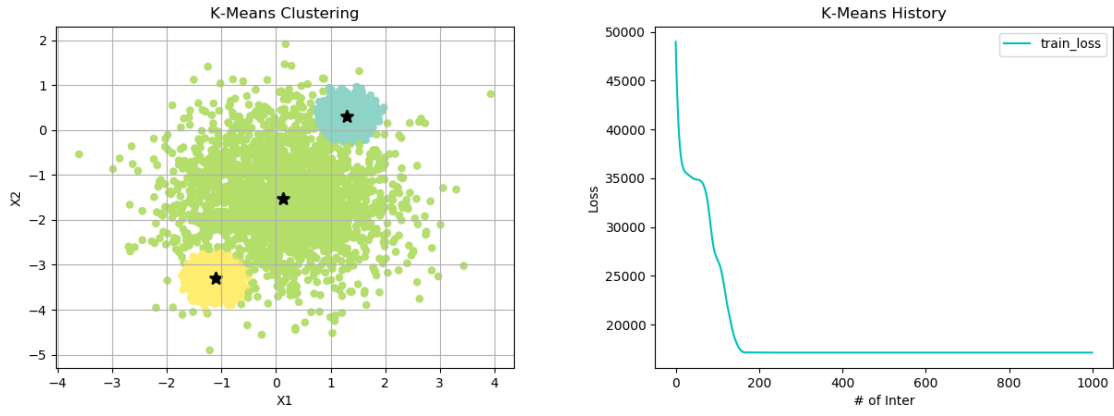


Figure 8: 2D cluster plot and Test data loss when $K = 3$.

Cluster	π	Mean	σ
1	-1.0983235	[1.299 0.308]	0.0388
2	-1.1029158	[-1.101 -3.307]	0.0391
3	-1.094615	[0.106 -1.527]	0.9872

Table 4: Best model parameters for test data set.

2.2.2

As shown in figure 10 below, the loss for validation data had decreased as k increases, however, the validation loss didn't have obvious decrease anymore from $k = 3$ onward, therefore, 3 is the best choice of k .

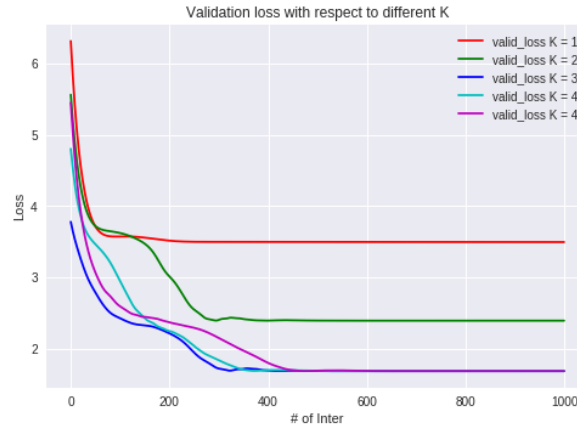


Figure 9: Validation loss with respect to different K .

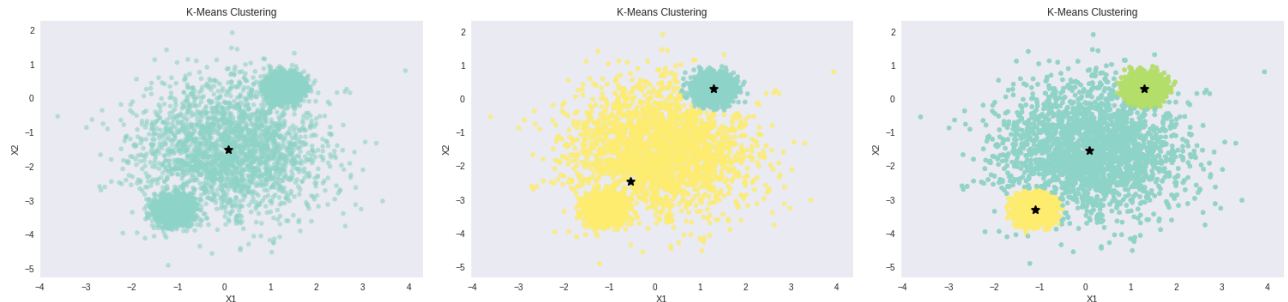


Figure 10: Validation loss with respect to different K (K from 1 to 3).

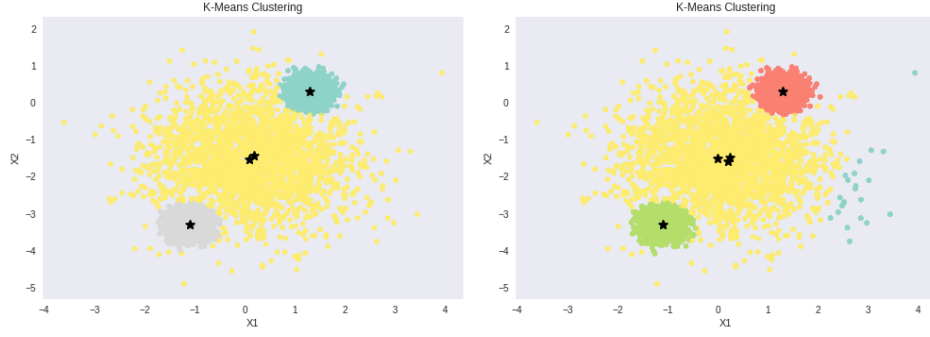


Figure 11: Validation loss with respect to different K(K from 4 to 5).

2.2.3

Figure 12 had showed the validation loss with different K by using GMM (at right hand-side), the loss values for all K values were similar after 1000 iterations of training, however, there was a large increase in the decay of loss values from $K = 5$ to $K = 10$, after that, the decay of loss values had not been changed much. Therefore, $K = 10$ is the reasonable amount of clusters for training using GMM.

In the other side, by looking at the loss values for different K values by using K-means, there was not large decrease after $K=10$, therefore, the reasonable amount of clusters is also 10. In addition, K-means didn't performed as well as GMM with high dimensional data.

Cluster K	5	10	15	20	25	30
Final training loss using GMM	50525.19	44320.33	44496.09	44057.45	44752.05	44147.00
Final training loss using K-means	247141.92	143545.07	143544.89	143545.48	143544.95	143544.96

Table 5: Final training loss with different K, *if_valid = True*.

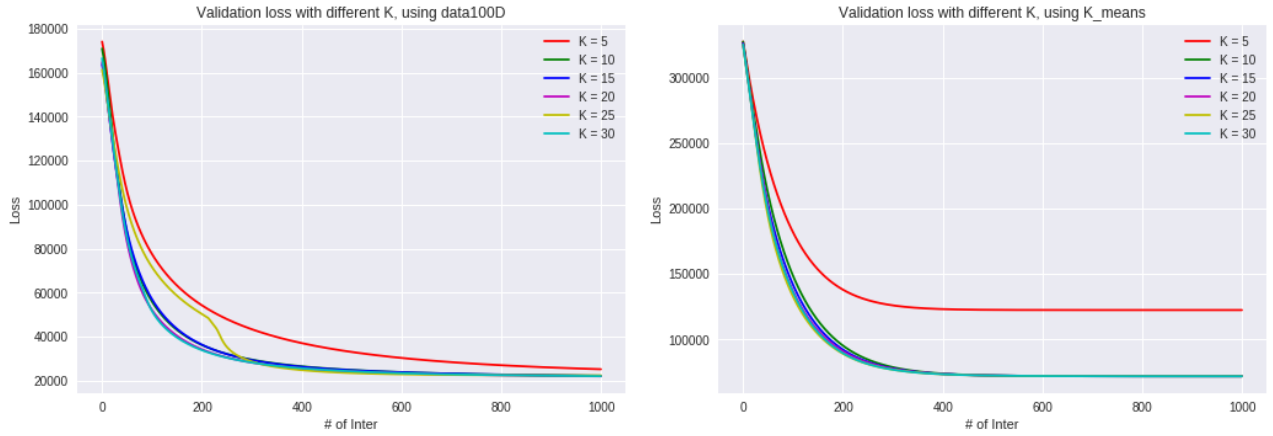


Figure 12: Validation loss, using GMM andK -means.